



FLASHLIGHT Workshop @ FCCM 2022 – May 18, 2022

Designing memory architectures with high-level synthesis: What could possibly go wrong?

Christian Pilato

Assistant Professor

christian.pilato@polimi.it

The EVEREST Project





EVEREST Approach

Big data applications with heterogeneous data sources

Three use cases





What are the relevant requirements for data, languages and applications?

How to design data-driven policies for computation, communication, and storage?

How to create FPGA accelerators and associated binaries?

How to manage the system at runtime?

How to evaluate the results?

How to disseminate and exploit the results?



Open-source framework to support the optimization of selected workflow tasks

FPGA-based architectures to accelerate selected kernels









EVEREST Use Cases





★ Improve the overall performance of traffic simulation

Accelerated computationally-intensive kernels

Machine-learning kernels



4

+

EVEREST Target System

cloudFPGA



- Disaggregated FPGAs directly attached to the network (64 FPGA instances)
- Low latency and high bandwidth system
- Separation between **Shell** and **Role** modules
- **cFDK framework** for system generation

FPGA-Accelerated HPC Cluster

- Cluster of **PCIe-attached FPGAs** (Alveo) with HBM architecture (up to 460 GB/s per board)
- Xilinx Vitis framework for HLS and system integration
- Support for the integration of **custom HDL**

CPU Reference System

CPU-based infrastructure to execute end-to-end workflows, manage storage, and data transfers
Extended to support the offloading of tasks to FPGA servers



Exploit spatial parallelism

High memory bandwidth

Different nodes to better match applications

Data-intensive (memory-bound) applications

Seamless support for multiple nodes



Limited **FPGA resources** (esp. memories)



EVEREST System Development Kit (SDK)



Different **input flows** starting from different **input languages**

Support for multiple target boards





Collection of interoperable and open-source tools to create hardware/software systems that can adapt to the target system, the application workflow, and the data characteristics

- Compilation framework based on **MLIR** to unify the input languages
- **High-level synthesis** and hardware generation flow to automatically create optimized architectures
- Creation of hardware and software variants to match architecture features
- Solution of state-of-the-art frameworks and commercial toolchains for FPGA synthesis



(and more...)

Ltvm

MIIR



The Case of Computational Fluid Dynamics

Numerical simulations are becoming more and more popular for many applications

- Computational Fluid Dynamics (CFD) is a representative application that requires to solve partial differential equations
- Kernel is the **Inverse Helmholtz operator** (parametric with respect to polynomial degree *p*) "Helmholtz" for the friends

_										
1	var	input	: S		:	[11	11]			
2	var	input	D		:	[11	11 11]			
3	var	input	u		:	[11	11 11]			
4	var	outpu	it v	7	:	[11	11 11]			
5	var	t			:	[11	11 11]			
6	var	r			:	[11	11 11]			
7	t =	S # 5	5 #	S	#	u.	[[1 6]	[3 7]	[5	8]]
8	r =	D * t	5							
9	v =	S # 5	5 #	S	#	t.	[[0 6]	[2 7]	[4	8]]

Final result is obtained by "small" contributions on independent data

- CFD kernel is composed of three high-level tensor operators (two contractions and one Hadamard product) repeated millions of times – good for spatial parallelism
- Each operator requires $p^2 + 2 \cdot p^3$ (double) elements as input and produces p^3 (double) elements 21.74 KB + 10.40 KB per element when p = 11
- It requires additional six tensors (p³ elements) to store intermediate results additional 62.39 KB



Hardware HPC (Memory) Architectures

What do we mean with **memory architecture**?

Every hardware module that is responsible to provide data to the accelerator kernels

Additional issues:

- BRAM resources are limited
 - Helmholtz operator requires >94 KB of local data
 - If local storage is not optimized, the number of parallel kernels can be limited
- Application-specific details can be used to optimize the data transfers
 - In Helmholtz, one of the tensors is constant over all elements how to match these details with platform characteristics?
 - Better to transfer data for a "batch" of elements and then execute them in series
 - how many? again, limited storage

Results on HBM FPGA – Spoiler Alert!





Challenges for HPC Architectures (i)

We can identify common challenges to most of the FPGA-based HPC architectures (e.g., network-attached cloudFPGA or bus-attached Alveo)

Challenge 1: Input languages and frameworks

• Application designers are usually not FPGA experts and may use high-level framework that are not supported by current HLS tools – how to talk with them?

Challenge 2: CPU-Host Communication Cost

• FPGA logic requires the data on the board, but data transfers can be much more expensive than kernel computation – how to execute more than one point?

Challenge 3: Read/Write Burst Transactions

• We need to determine the proper size of the transactions to get the maximum performance – how to reorganize the data transfers and get the parameters?



Challenges for HPC Architectures (ii)

Challenge 4: Full Bandwidth Utilization

- AXI interfaces may be large (e.g., 256 bits on the Alveo) how to leverage them?
- HBM architectures have many channels how to parallelize data transfers?

Challenge 5: Data Allocation

 Data must be placed in memory to maximize its utilization but also to enable efficient data transfers/computation –how to use custom data layouts?

Challenge 6: Synthesis-Related Issues

- FPGA devices are large but still not sufficient for hosting many kernels how to trade-off optimizations and parallel instances?
- FPGA (or architectures) may be different how to separate platform-agnostic and platform-dependent parts?
- FPGA logic architectures are complex and may introduce performance degradation – how to "guide" the synthesis process?



Hardware Compilation Flow





From DSL to Bitstream – Focus on Memory



POLITECNICO

MILANO 1863

PLM Customization for Heterogeneous SoCs

High-Level Synthesis (HLS) to create the accelerator logic

• Definition of memory-related parameters (e.g. number of process interfaces)

Generation of **specialized PLMs**

- Technology-related optimizations
- Possibility of system-level optimizations across accelerators





Reuse What is not Used

Generally, we use one **PLM unit** (possibly composed of many banks) for each **data structure** (array)

Reuse the same memory IPs for several data structures

"Two data structures are compatible if they can be allocated to the same PLM unit (memory IPs)"

<u>A common case</u>: accelerator parts never executed at the same time

- Possible only at system-level, when integrating the components
- Optimizations of accelerator logic and memory subsystem are independent



Memory Compatibility Graph (MCG)

Graph to represent the possibilities for optimizing the data structures

- Each node represents a data structure to be allocated, annotated with its data footprint (after data allocation)
- Each edge represents compatibility between the two data structures
- Can be automatically extracted from the MLIR-based compiler flow
 - Variant exploration to achieve the "best solutions"



- a) Address-space compatibility: the data structures are compatible and can use the same memory IPs
- b) Memory-interface compatibility: the ports are never accessed at the same time and the data structures can stay in the same memory IP



Preliminary Evaluation

- Xilinx Zynq UltraScale+ MPSoC ZCU106 board
 - CFD simulation of 50,000 elements

Memory sharing allows us to fit more kernels



MILANO 1863



Next Step: Let's Put Memory First

We are building an MLIR **compilation flow** for **automatic memory specialization**:

- MLIR Input DSL description of the system functionality
- Data Organization Determine which data resides off chip (also based on user/compiler annotations)
- **Layout** Reorganize communication to exploit local memories and perform efficient parallel computation
- Communication Configure prefetcher to hide transfer latency
- Local Partitioning Determine multi-bank PLM architecture (Mnemosyne)
- **HLS** Generate computation part (interfacing with existing HLS tools, e.g., open-source Vitis HLS frontend)
- HDL Output Automated code generation and system-level integration based on the target platform



S. Soldavini and C. Pilato. "Compiler Infrastructure for Specializing Domain-Specific Memory Templates" LATTE'21



Olympus – Automated System-Level Integration

We are developing a complete hardware architecture generation flow based on **MLIR description** of the **system functionality** Possibility to use several HLS

 \equiv **Platform-specific description** Kernel +Memorv Kernel Details Architecture (v) (json) Kernel Kernel IP HBM-based Xilinx Alveo Create IF \equiv Generation Black Box Kernel IBM CloudFPGA Description Platform CU \equiv Wrapper (C++ Olympus Backend Platform Platform Details Specific Opt 0101 001101 010011 \equiv HLS/Impl System Bitstream Host code generation Config Graph Graph CU Gen \equiv \equiv \equiv Extraction Transform Kernel Wrapper \equiv Based on platform libraries 0101 001101 010011 Compiler HW System General Platform Optimized -bound Description graph Binary for the specific target C++ CU (MLIR) graph Host Gen \equiv Host C++ Platform

POLITECNICO MILANO 1863 Host C++ Lib

tools/HDL generators

Conclusions

Data management optimizations are becoming the key for the creation of **efficient FPGA architectures** (... more than pure kernel optimizations)

HLS is now used not only to create accelerator kernels but also to generate the **system-level architecture**

• **Portable solutions** across multiple target platforms

Novel **HBM architectures** offer high bandwidth (that's why they are called *high-bandwidth memory* architectures... ③) but their design is complex:

- Necessary to match application requirements and technology characteristics
- We propose an MLIR-based compilation flow that directly interfaces with commercial HLS tools



Work done in collaboration with Stephanie Soldavini (Politecnico di Milano), Mattia Tibaldi (Politecnico di Milano), Jeronimo Castrillon (TU Dresden), Karl F. A. Friebel (TU Dresden), and Gerald Hempel (TU Dresden)

Thank you!

Christian Pilato, christian.pilato@polimi.it



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 957269